

Introduction:

This document introduces embedding techniques and vector databases as essential components for enhancing the relevance and performance of **Large Language Models (LLMs)**. By integrating external knowledge sources into LLM workflows, teams can build smarter, context-aware systems such as chatbots that provide more accurate, domain-specific responses.

Why Embeddings + Vector Databases?

While LLMs like GPT-3.5 turbo are powerful, they have critical limitations:

- Lack up-to-date information or custom domain knowledge.
- Tend to “hallucinate” when they don’t know something.
- Cannot natively “remember” large documents.

To address these challenges, embeddings are used. These are high-dimensional vector representations of text that capture semantic meaning. In this vector space, similar concepts are placed close together, allowing for comparison of meaning using mathematical similarity.

This approach forms the core of a strategy known as **Retrieval-Augmented Generation (RAG)** where relevant information is retrieved from an external vector database and added to the prompt before querying the LLM.

Embedding Techniques Overview:

This section introduces three widely used embedding techniques **Word2Vec**, **GloVe**, and **BERT** ; each representing a different generation of embedding methods

Word2Vec is one of the first methods to turn words into numbers (vectors) by looking at the words around them. It learns these patterns using a small neural network and creates one fixed vector for each word. It’s fast and efficient but can’t tell when a word has different meanings in different sentences.

GloVe works differently by counting how often words appear together across an entire text collection. This helps it understand broader relationships between words. Like Word2Vec, GloVe gives each word one fixed vector, so it also can’t adjust for different meanings based on context.

BERT is much more advanced. It looks at the full sentence and gives each word a vector that depends on its meaning in that sentence. So, the word “bank” will have a different

vector in “river bank” vs. “money bank.” This makes BERT great at understanding language, but it requires more computing power.

meaning, making it far more powerful for complex NLP tasks though more computationally expensive.

Technique	Type	Core Idea	Strengths	Limitations
Word2Vec	Static (Predictive)	Learns word meaning based on context using a shallow neural network (CBOW or Skip-gram).	<ul style="list-style-type: none">• Captures word analogies• Fast & lightweight• Pretrained models available	<ul style="list-style-type: none">• Same vector for every context• No sentence-level understanding• No support for OOV words
GloVe	Static (Count-based)	Creates word vectors by analyzing how often words appear together across a large collection of text.	<ul style="list-style-type: none">• Captures global semantics• Efficient lookup• Good analogy handling	<ul style="list-style-type: none">• Fixed vectors per word• No context awareness• Requires a large co-occurrence matrix
BERT	Contextual (Transformers)	Generates dynamic vectors for each word based on sentence context using bidirectional Transformers.	<ul style="list-style-type: none">• Context-aware embeddings• Excellent for sentence-level tasks• State-of-the-art accuracy	<ul style="list-style-type: none">• Slower & resource-intensive• Requires fine-tuning for semantic search• High memory usage

Embedding method of choice:

The project uses **Sentence-BERT**, a model that turns sentences into vectors based on their full meaning. Unlike older methods that give every word just one fixed vector, Sentence-BERT understands how a word’s meaning changes depending on the sentence. This helps the system find better matches between a user’s question and the stored document chunks. The project used the **all-MiniLM-L6-v2 model from Hugging Face**, which is fast, accurate, and works well with LangChain for building smart, responsive chatbots.

Vector Database for similarity search:

A **vector database** stores high-dimensional embeddings and enables **fast similarity search**. It powers **Retrieval-Augmented Generation (RAG)** by finding documents most similar in meaning to a user’s query.

Core Idea:

1. **User submits a question**
→ e.g. “What’s the return policy for damaged items?”
2. **Embed the question**
→ The question is converted into a numerical vector (embedding).
3. **Search vector database**
→ The system looks for chunks (pre-embedded document pieces) that are semantically similar.
4. **Retrieve top matches**
→ The most relevant chunks (e.g. policy snippets) are fetched as context.
5. **Construct enriched prompt**
→ These retrieved chunks are inserted into the prompt as additional context
6. **Send to LLM for generation**
→ The LLM uses this enriched prompt to generate a response that is **fact-based and grounded in the data**.

Instead of guessing or hallucinating, the chatbot retrieves facts from the data and uses it to give accurate, context-aware answers.

However, not all vector databases are created equal. Depending on the use case whether you're running a lightweight local prototype (chatbot) or deploying at scale in the cloud the choice of vector store can impact performance, flexibility, and cost.

The following section compares four widely used vector database solutions: **Facebook Artificial Intelligence Similarity Search (FAISS), Pinecone, Weaviate, and Qdrant**

Developer Need	Vector Database	Why is a good fit
Fast local development or prototyping	FAISS	Lightweight, in-process library with fast similarity search. Great for local code testing and demos.
Cloud-scale production deployment, zero setup	Pinecone	Fully managed SaaS. Scales to billions of vectors with minimal config or infrastructure effort.
Semantic + keyword (hybrid) search	Weaviate	Supports vector search. Ideal for nuanced enterprise queries.
Precise filtering and semantic search, easy setup	Qdrant	Rich metadata filtering + efficient Rust engine. Simple to deploy and self-host.

Vector Database of Choice: FAISS:

Justification:

- Easy to set up locally (no external infrastructure or hosting required)
- Fast and efficient for similarity search
- Well-integrated with **LangChain**, enabling quick experimentation with document embedding and retrieval

This made FAISS the **ideal choice for local development and testing** of the chatbot demo within the project's timeline and scope.

By combining **embedding models** and **vector databases**, teams can create **LLM applications** that are not only intelligent but also grounded in reliable, custom knowledge bridging the gap between general-purpose language models and domain-specific use cases.
